

JAN 07 2008

PATENT  
Atty. Dkt. No. YOR920010320US1**AMENDMENTS TO THE CLAIMS:**

This listing of claims will replace all prior versions, and listings, of claims in the application:

**LISTING OF CLAIMS**

1. (Previously Presented) A method, in a network comprising a primary server and a plurality of offload servers, for dynamic offloading of processing requests from said primary server to any one of said plurality of offload servers, the method comprising the steps of:

determining a load on said primary server;

if the load on said primary server is less than a first threshold, serving processing requests at said primary server; and

only if the load on said primary server exceeds said first threshold, then offloading at least a portion of said processing requests to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers.

2. (Original) The method of claim 1 wherein said load comprises bandwidth utilization and said first threshold is a network bandwidth utilization of said primary server.

3. (Original) The method of claim 1 wherein the said load comprises CPU utilization and said first threshold is a CPU utilization of said primary server.

4. (Previously Presented) The method of claim 1 wherein serving the processing requests at said primary server includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and  
offloading at least a portion of the processing requests to any one of said plurality

PATENT  
Atty. Dkt. No. YOR920010320US1

of offload servers includes serving a base page at said primary server in which the links for embedded objects point to any one of said plurality of offload servers.

5. (Previously Presented) The method of claim 1 wherein offloading at least a portion of the processing requests to any one of said plurality of offload servers includes routing an incoming Web request to a selected offload server.

6. (Original) The method of claim 1 and further including the step of, if the processing load on said primary server exceeds a second threshold, throttling at least one processing request.

7. (Previously Presented) The method of claim 6 wherein throttling at least one processing request includes returning a page to a user indicating that a server is overloaded.

8. (Original) The method of claim 6 wherein throttling at least one processing requests includes dropping the at least one processing request without returning any information to a user.

9. (Original) The method of claim 6 wherein throttling at least one processing request includes returning a page to a user indicating that a server is overloaded if said load exceeds said second threshold, and dropping said at least one processing request if said load exceeds a third threshold.

10. (Previously Presented) The method of claim 1 wherein the determination of which of said plurality of offload servers that at least one processing request is to be offloaded to is based on one or more of a group including a client identity, a client gateway (IP) address, a price of the offload service, or a current or previous load on the at least one offload server.

PATENT  
Atty. Dkt. No. YOR920010320US1

11. (Previously Presented) A network apparatus comprising a primary server and a plurality of offload servers connected by an IP-based network, for dynamic offloading of processing requests from said primary server to any one of said plurality of offload servers, the network apparatus comprising:

a load controller connected between said network and said primary server;

a memory connected to said load controller and including data and control instructions for operating said primary server to perform the steps of:

determining a load on said primary server;

if said load on said primary server is less than a first threshold, serving processing requests at said primary server; and

only if said load on said primary server exceeds said first threshold, then offloading at least a portion of said processing requests to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers.

12. (Original) The apparatus of claim 11 wherein said load comprises network bandwidth and said first threshold is a network bandwidth utilization of said primary server.

13. (Original) The apparatus of claim 11 wherein said load comprises CPU utilization and said first threshold is a CPU utilization of said primary server.

14. (Previously Presented) The apparatus of claim 11 wherein serving the processing requests at said primary server includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and

offloading at least a portion of the processing requests to any one of said plurality of offload servers includes serving a base page at said primary server in which the links for embedded objects point to any one of said plurality of offload servers.

PATENT  
Atty. Dkt. No. YOR920010320US1

15. (Previously Presented) The apparatus of claim 11 wherein offloading at least a portion of the processing requests to any one of said plurality of offload servers includes routing an incoming Web request to a selected offload server.

16. (Original) The apparatus of claim 11 and further including, if the processing load on said primary server exceeds a second threshold, throttling at least one processing request.

17. (Previously Presented) The apparatus of claim 16 wherein throttling at least one processing request includes returning a page to a user indicating that a server is overloaded.

18. (Original) The apparatus of claim 16 wherein throttling at least one processing requests includes dropping the at least one processing request without returning any information to a user.

19. (Original) The apparatus of claim 16 wherein throttling at least one processing request includes returning a page to a user indicating that a server is overloaded if the primary server load exceeds said second threshold, and dropping the at least one processing request if said primary server load exceeds a third threshold.

20. (Previously Presented) The apparatus of claim 11 wherein the determination of which of said plurality of offload servers that at least one processing request is to be offloaded to is based on one or more of a group including a client identity, a client gateway (IP) address, a price of the offload service, or a current or previous load on the at least one offload server.

21. (Previously Presented) A system, including an IP network comprising a primary server and a plurality of offload servers, for dynamic offloading of processing requests

from said primary server to any one of said plurality of offload servers, the system comprising:

means for determining a load on said primary server;

means for, if said load on said primary server is less than a first threshold, serving the processing requests at said primary server; and

means for, only if said load on said primary server exceeds said first threshold, offloading at least a portion of said processing requests to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers.

22. (Previously Presented) A program product including control instructions for controlling the operation of a computer, said program product operative with said control instructions to operate said computer in an IP-based network comprising a primary server and a plurality of offload servers to dynamically offload processing requests from said primary server to any one of said plurality of offload servers, the computer operative with said control instructions to perform the steps of:

determining a load on said primary server;

if said load on said primary server is less than a first threshold, serving the processing requests at said primary server; and

only if said load on said primary server exceeds said first threshold, then offloading at least a portion of said processing requests to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers.

23. (Previously Presented) A system for allocating processing requirements on a network between a primary server and a plurality of offload servers, comprising:

a load controller connected to said network for receiving processing requests from clients on said network and allocating said processing requests between said

**PATENT**

Atty. Dkt. No. YOR920010320US1

primary and offload servers;

a memory connected to said load controller and storing threshold data and control software for analyzing said threshold data and operating said load controller;

said load controller operative with the threshold data and control software to perform the steps of:

periodically evaluating said processing requests to determine a load on said primary server;

if said load exceeds a first threshold, for a predetermined period of time directing at least one processing request to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers; and

only if said load does not exceed said first threshold, directing said processing requests to said primary server.

24. (Original) The system of claim 23 wherein said load is network bandwidth and said first threshold is a measure of the network bandwidth utilization of the primary server.

25. (Original) The system of claim 23 wherein said load is CPU utilization and said first threshold is a measure of the CPU utilization of the primary server.

26. (Previously Presented) The system of claim 23 wherein directing said processing requests to said primary server further includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and

directing at least one processing request to any one of said plurality of offload servers further includes serving a base page at said primary server in which the links for embedded objects point to said any one of said plurality of offload servers.

27. (Previously Presented) The system of claim 23 wherein directing at least one

PATENT  
Atty. Dkt. No. YOR020010320US1

processing request to any one of said plurality of offload servers further includes routing an incoming Web request to a selected offload server.

28. (Previously Presented) The system of claim 23 and further including, if said load exceeds a second threshold, throttling at least one processing request by returning a page to a user indicating that a server is overloaded.

29. (Original) The system of claim 23 and further including, if said processing load exceeds a second threshold, dropping said at least one processing request without returning any information to a user.

30. (Original) The system of claim 23 and further including throttling at least one processing request by returning a page to a user indicating that said primary server is overloaded if said primary server load exceeds a second threshold, and dropping the at least one processing request if said primary server load exceeds a third threshold.

31. (Previously Presented) The system of claim 23 further including determining which of said plurality of offload servers that said at least one processing request is to be offloaded to is based on one or more of a group including a client identity, a client gateway (IP) address, a price of the offload service, or a current or previous load on the at least one offload server.

32. (Previously Presented) A method for allocating processing requirements on an IP network between a primary server and a plurality of offload servers, comprising:  
periodically evaluating processing requests to determine a load on said primary server;

if said load exceeds a first threshold, for a predetermined period of time directing at least one processing request to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled

PATENT  
Atty. Dkt. No. YOR920010320US1

by said plurality of offload servers; and

only if said processing load does not exceed said first threshold, directing said processing requests to said primary server.

33. (Original) The method of claim 32 wherein said load comprises network bandwidth and said first threshold is a measure of the network bandwidth utilization of said primary server.

34. (Original) The method of claim 32 wherein said load comprises CPU utilization and said first threshold is a measure of the CPU utilization of said primary server.

35. (Previously Presented) The method of claim 32 wherein directing said processing requests to said primary server further includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and

directing at least one processing request to any one of said plurality of offload servers further includes serving a base page at said primary server in which the links for embedded objects point to said any one of said plurality of offload servers.

36. (Previously Presented) The method of claim 32 wherein directing at least one processing request to any one of said plurality of offload servers further includes routing an incoming Web request to a selected offload server.

37. (Previously Presented) The method of claim 32 and further including, if said load exceeds a second threshold, throttling at least one processing request by returning a page to a user indicating that a server is overloaded.

38. (Original) The method of claim 32 and further including, if said load exceeds a second threshold, dropping the at least one processing request without returning any information to a user.



39. (Original) The method of claim 32 and further including throttling at least one processing request by returning a page to a user indicating that the primary server is overloaded if the primary server load exceeds a second threshold, and dropping the at least one processing request if the primary server load exceeds a third threshold.

40. (Previously Presented) The method of claim 32 further including determining which of said plurality of offload servers that said at least one processing request is to be offloaded to is based on one or more of a group including a client identity, a client gateway (IP) address, a price of the offload service, or a current or previous load on the at least one offload server.

41. (Previously Presented) A system for allocating processing requirements on an IP network between a primary server and a plurality of offload servers, comprising:

means for periodically evaluating processing requests to determine a load on said primary server;

means for, if said load exceeds a first threshold, for a predetermined period of time directing at least one processing request to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers; and

means for, only if said processing load does not exceed said first threshold, directing said processing requests to said primary server.

42. (Previously Presented) A program product including control instructions for controlling the operation of a computer, said program product operative with said control instructions to operate said computer in an IP-based network comprising a primary server and a plurality of offload servers to dynamically offload processing requests from said primary server to any one of said plurality of offload servers, the computer operative with said control instructions to perform the steps of:

periodically evaluating processing requests to determine a load on said primary

PATENT  
Atty. Dkt. No. YOR920010320US1

server;

if said load exceeds a first threshold, for a predetermined period of time directing at least one processing request to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers; and

only if said load does not exceed said first threshold, directing said processing requests to said primary server.